



**caBIG**

*cancer Biomedical  
Informatics Grid*



# Overview of ICR Workspace

Thomas L. Casavant, Univ of Iowa/Holden

SP-SLWG Face to Face meeting

November 7-8, 2004

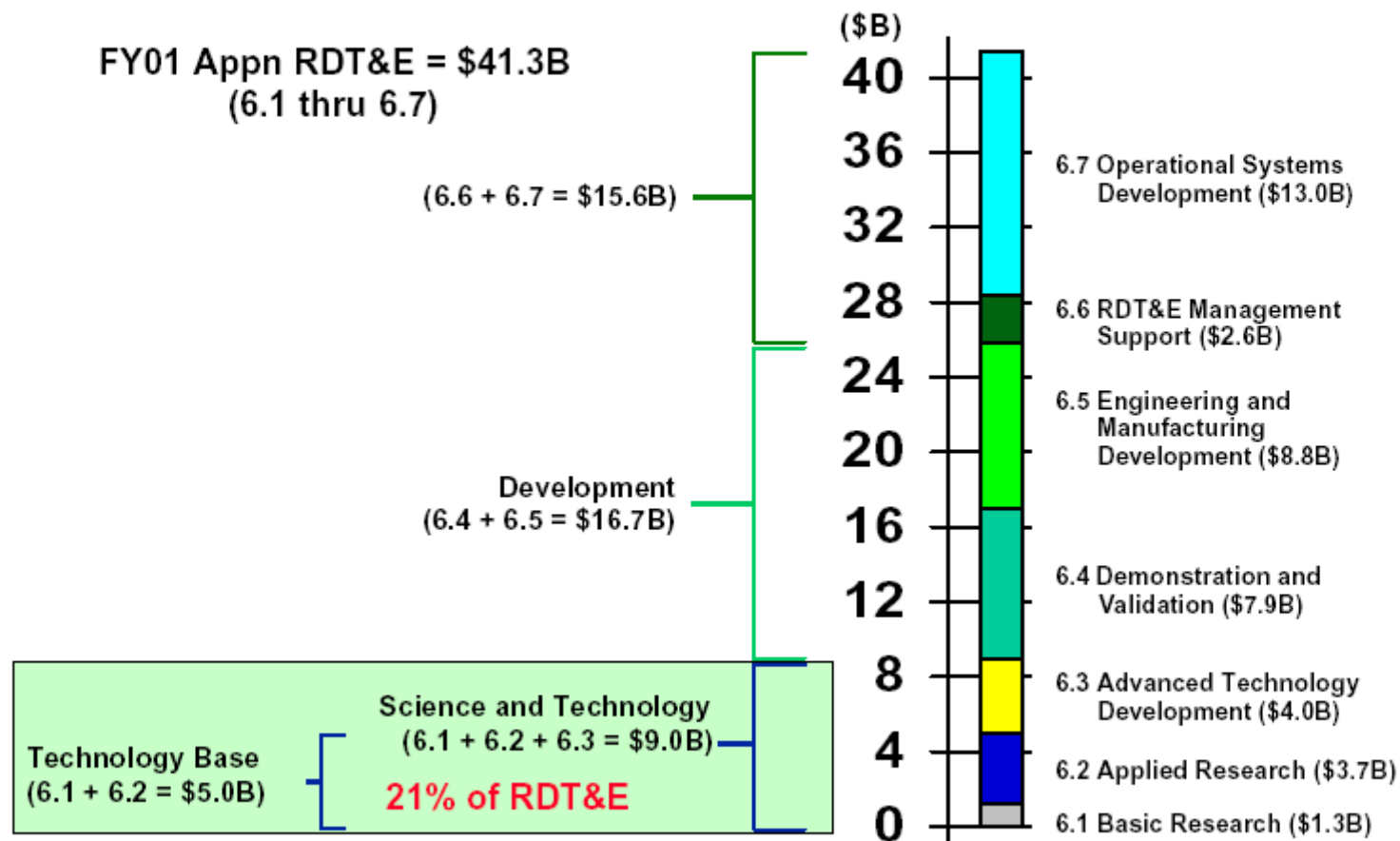
# ICR Workspace at a Glance

## - Overview::Definition

- ▶ The Integrative Cancer Research Workspace is producing modular and interoperable tools and interfaces that provide for integration between biomedical informatics applications and data. This will ultimately enable translational and integrative research by providing for the integration of clinical and basic research data. The Workspace is developing a software-engineered, well-documented and validated biomedical informatics toolset for use throughout the research community.
- ▶ DoD/DARPA-speak
  - I see this as a classic 6.1-6.2 charter, whereas CTMS is more 6.6



# FY01 Appropriated RDT&E



# ICR Workspace at a Glance

## - Overview::Composition

- ▶ Workspace Lead: Juli Klemm
- ▶ SP-SLWG Liaison – Michael Ochs- Fox Chase
- ▶ 6 ICR SIGs (8 proposed in February)
- ▶ ICR Membership
  - 152 on email alias
  - 69 attended 1<sup>st</sup> F2F meeting in Bethesda, August 25-26, 2004
    - 32 Institutions, plus
    - 14 from BAH and/or NCICB/SAIC
    - 9 from NCICB
    - 1 Oracle, 1 Alpha Gamma Technologies
- ▶ 21 ICR Development Projects
  - 19 Silver
  - 2 Gold reference implementations
- ▶ July/August caBIG program update highlighted ICR

# ICR Workspace at a Glance

## - Overview::Funded Participants

16 Developers:	7 Adopters:	4 Participants:
Columbia	Wistar	Vanderbilt -- Ingram
UNC Lineberger	New York University	U. of Michigan
UC San Francisco	Penn – Abrahamson	Prentis – Karmanos
Georgetown -- Lombardi	Memorial Sloan Kettering	Northwestern -- Lurie
Burnham Institute	Georgetown	
Wash U. -- Siteman	Oregon Health	
NCI – Center for Cancer Reseach	U. of South Florida -- Moffitt	
Cold Spring Harbor		
U. of Chicago		
Thomas Jefferson – Kimmel		
Memorial Sloan Kettering		
Fox Chase		
Dartmouth		
Duke		
University of Iowa -- Holden		
MIT/Broad		

# ICR Workspace at a Glance

## - Overview::Liaisons

Tissue Banks and Pathology  
Clinical Trials Management  
Architecture  
Vocabularies and Common Data Elements  
Data Sharing and Intellectual Capital  
Training  
Strategic Planning

Wash U Mark Watson  
Penn David Fenstermacher  
Duke Patrick McConnell  
Vanderbilt Mary Edgerton  
UI/Holden Tom Casavant/Terry Braun  
Institute for Cancer Prevention Edith Zang  
Fox Chase Michael Ochs

# ICR Workspace at a Glance

## - Overview::Composition::6 SIGs

SIG	(#projects, #members, #on a typical conference call)
1. Data Analysis and Statistical Tools	(5, 52, 15-20)
– Lead: Chris Kingsley, UCSF	
2. Genome Annotation	(6, 41, 10-15)
– Lead: Craig Street, Penn	
3. Microarray Repositories	(3, 47, 15)
– Lead: Julie Zhu, Northwestern	
4. Pathways	(3, 35, 10)
– Lead: Shannon McWeeney, OHSU	
5. Proteomics	(3, 43, 10)
– Lead: Sinoula Apostolou, Fox Chase	
6. Translational	(1, 40, 10)
– Lead: Terry Braun, U Iowa/Holden	

# ICR Workspace at a Glance

## - Overview::Brief Chronology Since Kickoff

February	<ul style="list-style-type: none"> <li>- Kickoff Meeting</li> <li>- Special Interest Groups Defined</li> </ul>
April	<ul style="list-style-type: none"> <li>- First Workspace Teleconference</li> <li>- Juli Klemm starts as Workspace Lead</li> <li>- Begin conversations with Centers about projects</li> </ul>
May	<ul style="list-style-type: none"> <li>- Special Interest Group meetings begin</li> <li>- Continue project definitions</li> <li>- Begin Developer-Adopter pairings</li> <li>- Liaisons to other Workspaces identified</li> <li>- SOW template established</li> </ul>
June	<ul style="list-style-type: none"> <li>- Begin drafting SOWs</li> <li>- Project presentations begin in SIGs</li> <li>- Finalize Developer-Adopter pairings</li> </ul>
July	<ul style="list-style-type: none"> <li>- ICR featured in caBIG Program Update</li> <li>- Begin identification of candidate grid reference implementations</li> <li>- Continue drafting SOWs</li> </ul>
August	<ul style="list-style-type: none"> <li>- ICR face-to-face meeting</li> <li>- First Workspace Participant task orders issued</li> </ul>
September	<ul style="list-style-type: none"> <li>- Begin cost negotiations with Centers</li> <li>- First RFP issued (Duke)</li> <li>- Cancer Center leads for SIGs established</li> </ul>
October	<ul style="list-style-type: none"> <li>- Six more RFPs issued</li> <li>- First projects started: NCI-CCR's GoMiner and NCI-60 projects</li> <li>- VCDE liaisons established for SIGs</li> <li>- Developer/Adopter teams begin teleconferences, scheduling face-to-face meetings</li> <li>- SIGs begin discussion of exchange standards; Gene CDE focus group established</li> <li>- Begin UML-based CDE creation training</li> </ul>



# ICR Workspace at a Glance

## - Overview::Meeting Schedule

- ▶ General Workspace Meetings
  - 2nd Wednesday of each month, 2:00pm - 3:00pm
- ▶ Data Analysis & Statistical Tools SIG
  - 1st Friday of each month, 2:00pm - 3:00pm
- ▶ Genome Annotation SIG
  - 1st Thursday of each month, 3:00pm - 4:00pm
- ▶ Informatics for Proteomics SIG
  - 2nd Monday of each month, 2:00pm - 3:00pm
- ▶ Microarray Repositories SIG
  - 1st Wednesday of each month, 2:00pm - 3:00pm
- ▶ Pathways Tools SIG
  - 1st Tuesday of each month, 1:00pm - 2:00pm
- ▶ Translational Tools SIG
  - 1st Monday of each month, 1:00pm - 2:00pm

# **SIG 1:**

## **Data Analysis and Statistical Tools**

**(5, 52, 15-20)**

The mission of the Data Analysis Special Interest Group of the ICR Workspace is to serve the needs of key categories of end users - experimentalists and data analysts - by provision of interoperable tools and associated standards, documentation, and training.

- Data analysis was identified as one of the most substantial needs across the Cancer Centers in the formal assessment that preceded initiation of the caBIG development workspaces.
- Much of the need stems from the increased complexity and volume of data sets resulting from high-throughput measurement technologies.

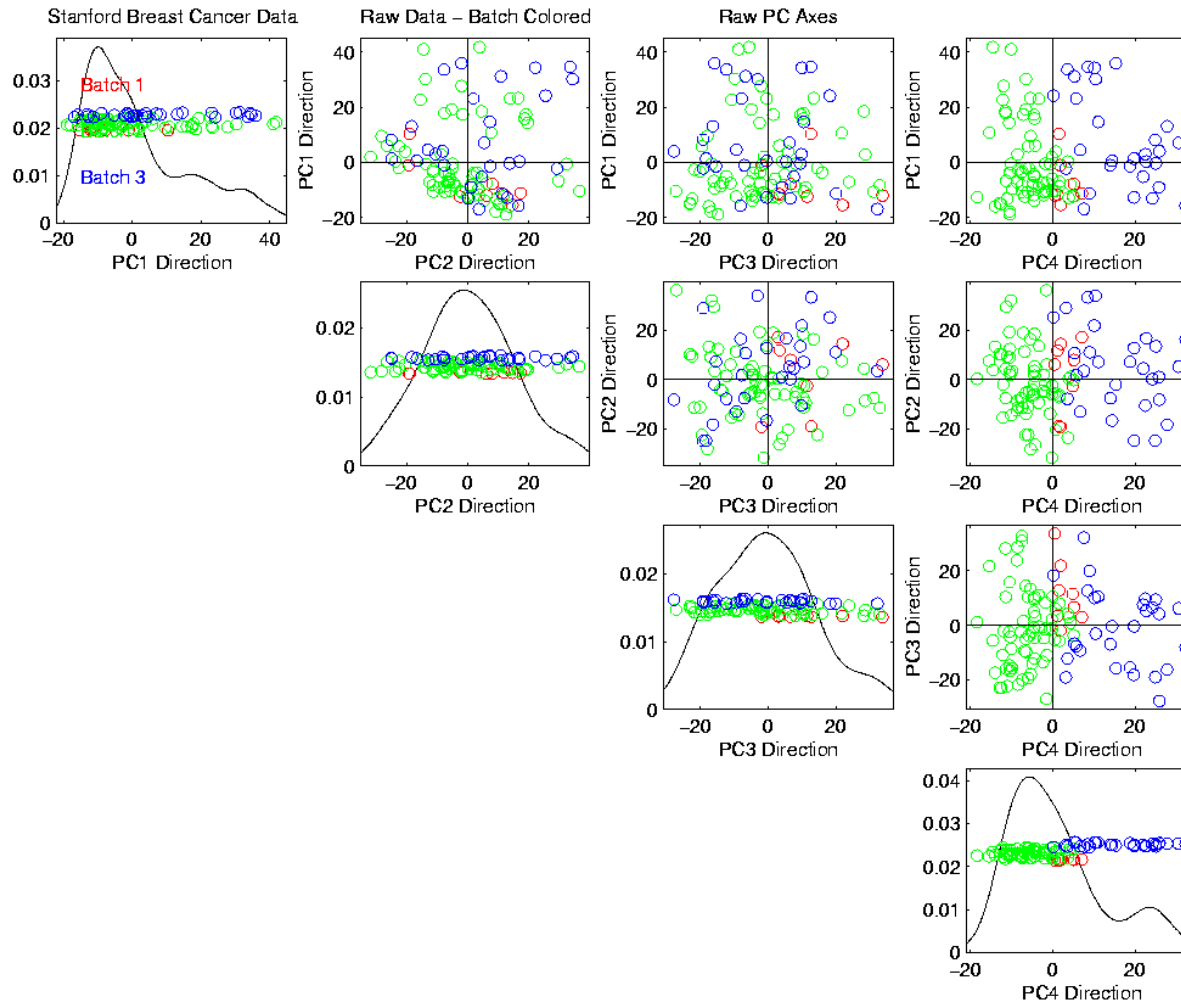
## Distance-Weighted Discrimination (DWD)

<b>Developer Center -- POC</b>	UNC Lineberger – Steve Marron
<b>Adopter Center -- POC</b>	Wistar – Louise Showe
<b>caBIG Compliance Level to be Achieved</b>	Silver
<b>Project duration</b>	12 months

# DWD Project Objectives

- ▶ Create a Risk Management Matrix for the project
- ▶ In collaboration with Adopter, create a use case document for DWD
- ▶ Develop Functional Requirements and Design Specification document, in collaboration with the Adopter Center and the cross-cutting Workspaces
- ▶ Document a Test Approach that ensures requirements are met
- ▶ Perform preliminary evaluation of Adopter's data
- ▶ Code caBIG-compliant DWD
- ▶ Create visual diagnostics for DWD
- ▶ Deploy caBIG-compliant DWD to Adopter site
- ▶ Execute on Test Approach

# Source Batch Adj: Batch Colors



# GenePattern

<b>Developer Center -- POC</b>	MIT/Broad – Jill Mesirov
<b>Adopter Center -- POC</b>	NYU -- Judith Goldberg
<b>caBIG Compliance Level to be Achieved</b>	Silver
<b>Project duration</b>	12 months

# GenePattern

## GenePattern

GenePattern Home

Download

FAQ

Tutorial

Algorithms

Data Sets

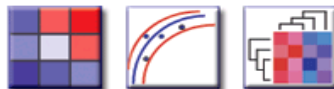
Mailing List

SEARCH

**GenePattern** is a flexible analysis platform developed to support multidisciplinary biomedical research. GenePattern puts the power of sophisticated computational methods into the hands of non-programming users. It also provides an environment for rapid development and deployment of new analytic techniques.

10.8.2004 A new version of GenePattern has been released. GenePattern 1.2.1 is available [here](#)

### Graphical Environment



An **intuitive user interface** provides extensive support for users at all levels of computational sophistication.

### Pipeline Environment



Users can **chain tasks together** to create, encapsulate, reproduce, and share methodologies.

### Analysis Tool Repository



**Add your own modules** to the GenePattern environment without writing extra code, or choose from a **comprehensive repository of analysis modules**.

### Programming Language Environment



Computational biologists and software developers have **programmatic access** to all GenePattern modules from any of several programming languages.

GenePattern is funded by a grant from the [NIH](#).  
Copyright © 2004 Broad Institute of MIT and Harvard

# Magellan

<b>Developer Center -- POC</b>	UCSF – Ajay Jain
<b>Adopter Center -- POC</b>	Penn -- David Fenstermacher
<b>caBIG Compliance Level to be Achieved</b>	Silver
<b>Project duration</b>	12 months



# Magellan Project Objectives

- ▶ Develop Functional Requirements and Design Specification documents, in collaboration with the Adopter Center(s) and the cross-cutting Workspaces
- ▶ Create a Risk Management Matrix for the project
- ▶ Document a Test Approach that ensures requirements are met
- ▶ Write code to achieve the following milestones:
  - Interoperability between Magellan and caBIO
  - Interoperability between Magellan and caArray
- ▶ Execute on Test Approach
- ▶ Deploy system to Adopter site

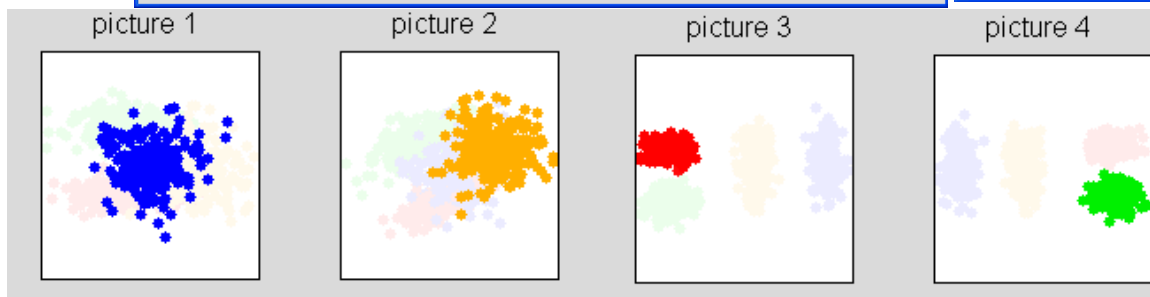
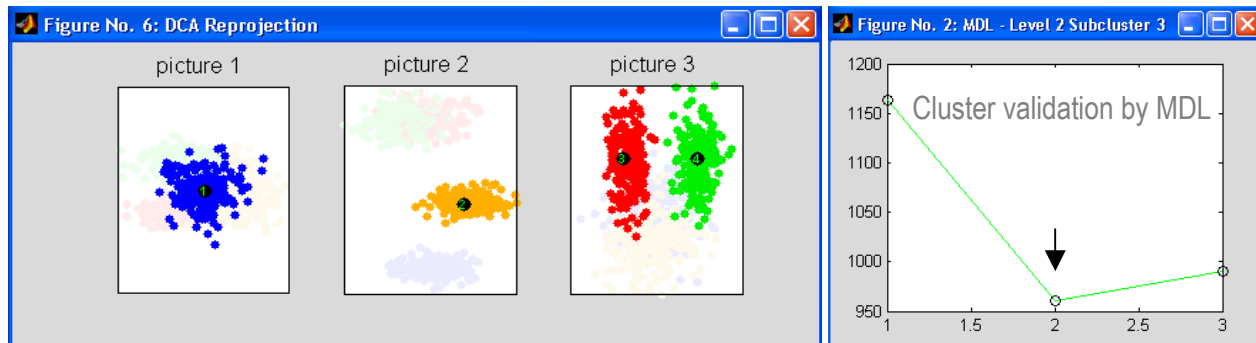
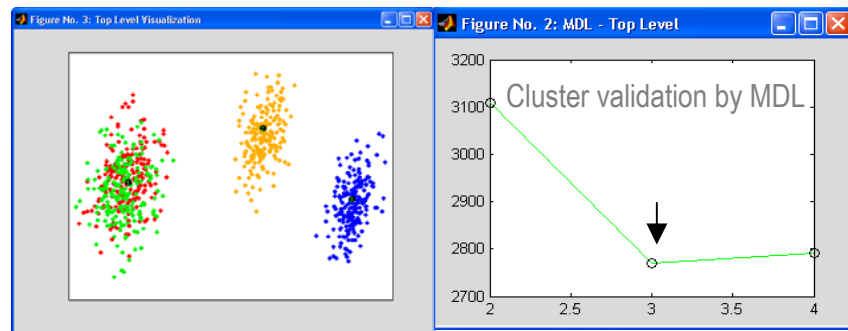
## Visual Statistical Data Analyzer -- VISDA

<b>Developer Center -- POC</b>	Georgetown – Joseph Wang
<b>Adopter Center -- POC</b>	Wistar – Louise Showe
<b>caBIG Compliance Level to be Achieved</b>	Silver
<b>Project duration</b>	12 months

# VISDA Project Objectives

- ▶ Develop a Functional Requirements and Design Specification document, in collaboration with the Adopter Center(s) and the cross-cutting Workspaces
- ▶ Create a Risk Management Matrix for the project
- ▶ Document a Test Approach that ensures requirements are met
- ▶ Write code for the following functionality:
  - Cluster modeling using hierarchical mixture modeling
  - Dimension reduction by principle component analysis (PCA), discriminatory component analysis (DCA) and project pursuit method (PPM)
  - Cluster formation by soft data clustering using expectation-maximization (EM) algorithm
  - Cluster validation by minimum description length (MDL) criterion
  - Cluster visualization by hierarchical cluster display
  - Graphical user interface (GUI) for VISDA set-up, data analysis, and data visualization
- ▶ Execute on Test Approach

# “Simulation with Truth”



## **SIG 2: Genome Annotation (6, 41, 10-15)**

The mission of the Genome Annotation SIG is to provide data and tools that will greatly enhance the cancer research community's access to high quality, comprehensive gene annotations. Having standardized access to these data sources will support a consistent view of all available gene information. This will be achieved by adapting existing software that meets the needs of the user community to comply with caBIG and by creating new software and tools.

## Cancer Molecular Pages

<b>Developer Center -- POC</b>	Burnham – Kutbuddin Doctor
<b>Adopter Center -- POC</b>	TBD
<b>caBIG Compliance Level to be Achieved</b>	Silver
<b>Project duration</b>	12 months

# Cancer Molecular Pages Project Objectives

- ▶ Develop use cases for Cancer Molecular Pages
- ▶ Develop Functional Requirements and Design Specification documents, in collaboration with the Adopter Center(s) and the cross-cutting Workspaces
- ▶ Create a Risk Management Matrix for the project
- ▶ Document a Test Approach that ensures requirements are met
- ▶ Create a prototype of the Cancer Molecular Pages application and deploy at Adopter site
- ▶ Execute on Test Approach
- ▶ Install caArray and migrate legacy data to the system

## FunctionExpress

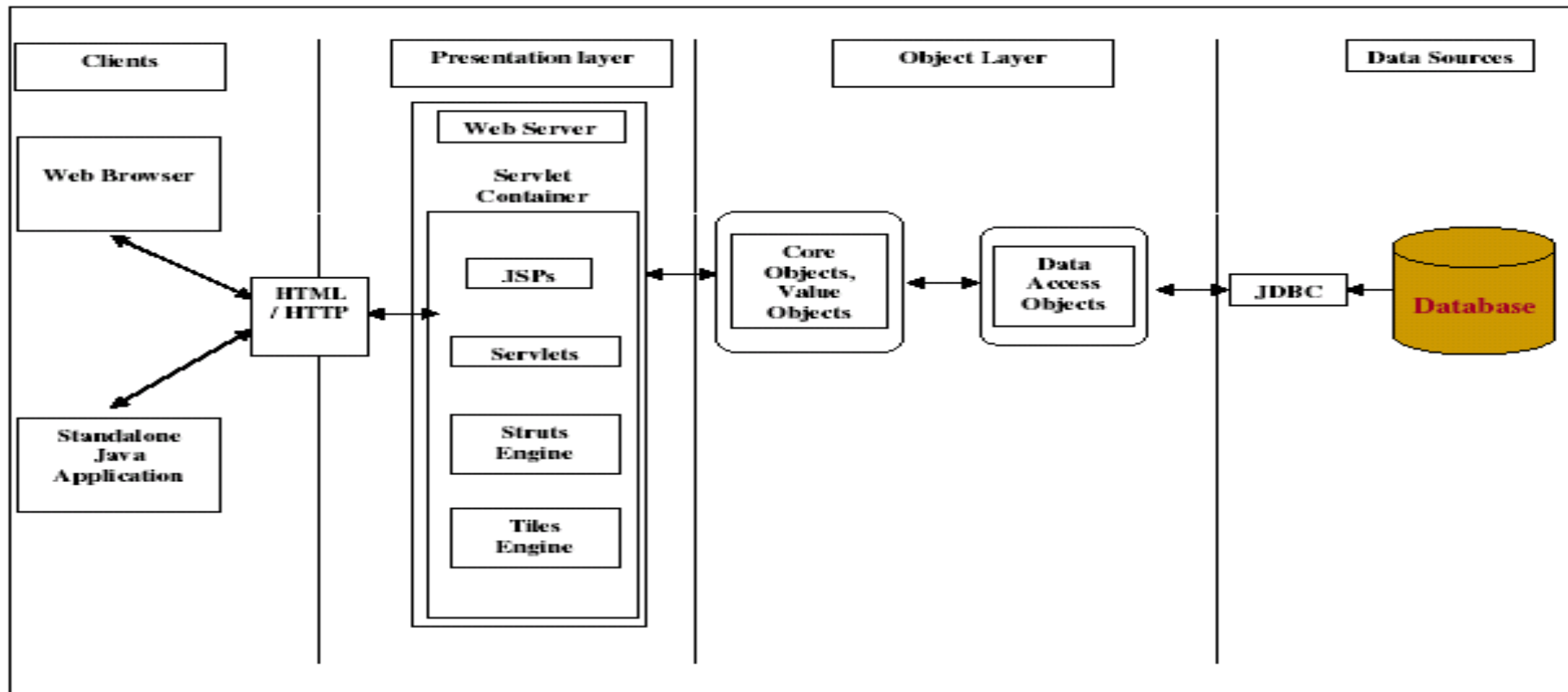
<b>Developer Center -- POC</b>	Wash U – Rakesh Nagarajan
<b>Adopter Center -- POC</b>	Wistar – Harold Riethman
<b>caBIG Compliance Level to be Achieved</b>	Silver
<b>Project duration</b>	12 months



# FunctionExpress Project Objectives

- ▶ Develop a Use Case document, in collaboration with the Adopter Center(s)
- ▶ Develop a Functional Requirements collaboration with the Adopter Center(s)
- ▶ Develop a Design Specification document and UML class diagram of the system.
- ▶ Create a Risk Management Matrix for the project
- ▶ Document a Test Approach that ensures requirements are met
- ▶ Write code to achieve the following goals:
  - Implementation of the object and data models for the web application
  - Implement web application for annotation-based gene search and for display of gene literature networks .
  - Import microarray data from external sources, including caArray
  - Import annotation data from external sources
- ▶ Execute on Test Approach
- ▶ Deploy beta system to Adopter site

# Architecture Diagram of FE



**High Level Architecture Diagram of Web Application**

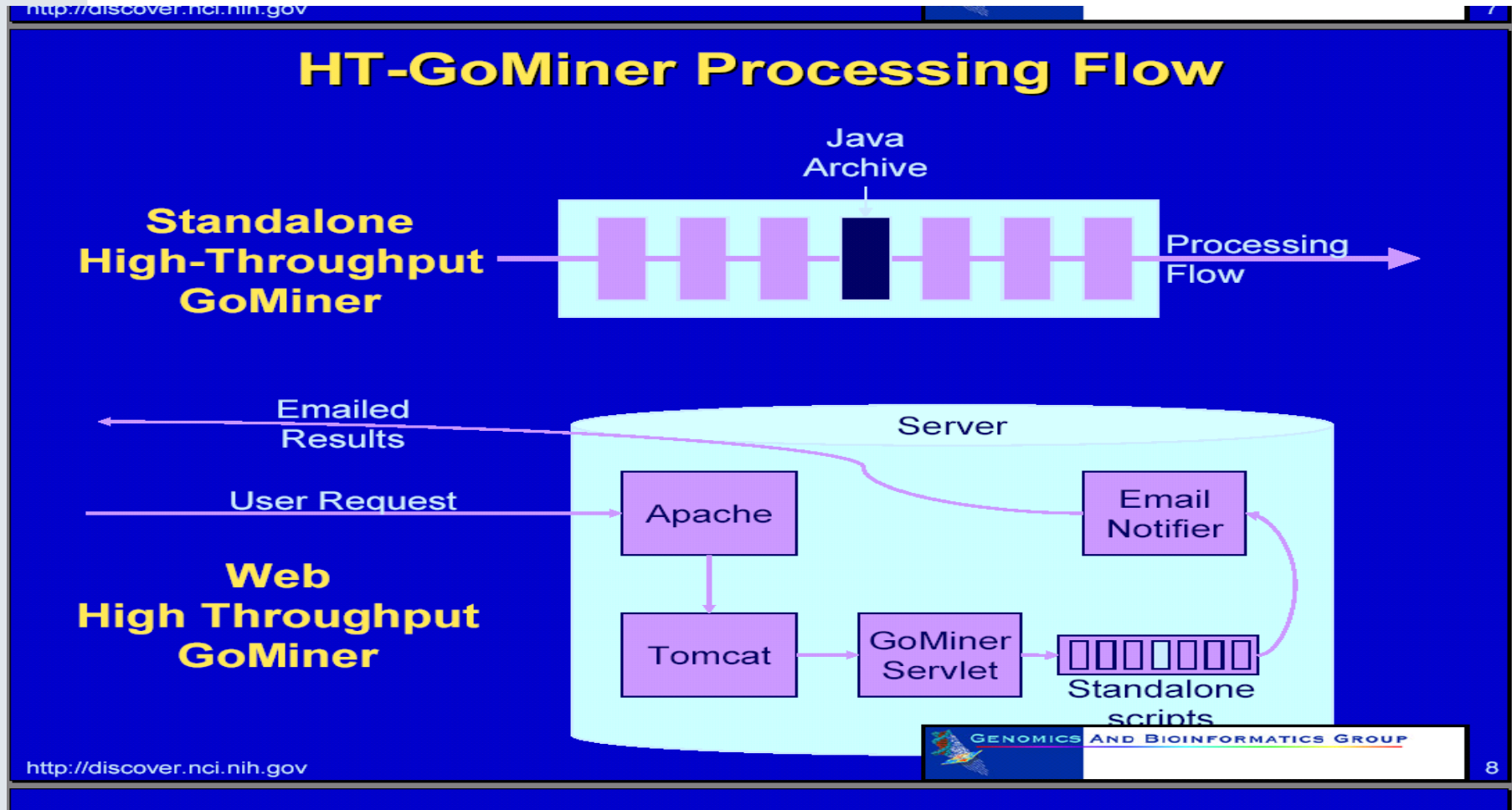
## GoMiner

<b>Developer Center -- POC</b>	NCI-CCR – David Kane
<b>Adopter Center -- POC</b>	Wistar – Harold Riethman
<b>caBIG Compliance Level to be Achieved</b>	Silver
<b>Project duration</b>	8 months

# GoMiner Project Objectives

- ▶ Create a use case document for GoMiner, including the use of BioCarta for classification
- ▶ Develop Functional Requirements and Design Specification documents, in collaboration with the Adopter Center(s) and the cross-cutting Workspaces.
- ▶ Create a Risk Management Matrix for the project
- ▶ Document a Test Approach that ensures requirements are met
- ▶ Write code to achieve the following milestones:
  - Implement an XML-based output format with explicit meta-data descriptors
  - Implement a mechanism to retrieve data from caArray with a mechanism to specify a desired data set
  - Develop a web services API for GoMiner
- ▶ Execute on Test Approach
- ▶ Package and release the GoMiner source code under an Open Source license

# High-Throughput GoMiner Processing Flow



## HapMap, Vertebrate Promoter DB

<b>Developer Center -- POC</b>	Cold Spring Harbor – Brian Gilman
<b>Adopter Center -- POC</b>	Wistar – Harold Riethman MSKCC- Alex Lash
<b>caBIG Compliance Level to be Achieved</b>	Silver
<b>Project duration</b>	9 months

# HapMap, Vertebrate Promoter DB Project Objectives

- ▶ Develop a Functional Requirements and Design Specification document, in collaboration with the Adopter Center(s) and the cross-cutting Workspaces
- ▶ Create a Risk Management Matrix for the project
- ▶ Document a Test Approach that ensures requirements are met
- ▶ Write code to achieve the following milestones:
  - Establish use cases for the two databases
  - Serve HapMap and the VPD by the DAS database
  - Import HapMap/VPD Data into the caBIO database
  - Augment caBIO system to serve Hapmap and VPD data
  - Validate import of HapMap and VPD data to caBIO object model
- ▶ Execute on Test Approach

## Protein Information Resource (PIR)

<b>Developer Center -- POC</b>	Georgetown -- Cathy Wu
<b>Adopter Center -- POC</b>	Penn – David Fenstermacher
<b>caBIG Compliance Level to be Achieved</b>	Gold
<b>Project duration</b>	7 months (for first phase of project)



# PIR Project Objectives

- ▶ In collaboration with the Adopter center, develop use cases for grid-enabled PIR
- ▶ Develop Functional Requirements and Design Specification documents, in collaboration with the Adopter Center(s) and the cross-cutting Workspaces
- ▶ Create a Risk Management Matrix for the project
- ▶ Document a Test Approach that ensures requirements are met
- ▶ Provide the following to support a PIR middleware layer:
  - Database to object mapping description
  - Produce a web services layer to access data
  - Publish description of service in WSDL with XML schema that defines input and output objects
- ▶ Execute on Test Approach

# UniProt Report

## PIR View

UniProt Entry: **P00439**

## ENTRY INFORMATION

ENTRY NAME	<a href="#">PH4H_HUMAN</a>
ACCESSION NUMBERS	P00439; Q16717; Q8TC14
CREATED	Release 01, 21-JUL-1986
SEQUENCE UPDATE	Release 01, 21-JUL-1986
ANNOTATION UPDATE	Release 44, 05-JUL-2004

## NAME AND ORIGIN OF THE PROTEIN

PROTEIN NAME	Phenylalanine-4-hydroxylase
DESCRIPTION	(EC 1.14.16.1) . PAH; Phe-4-mono-oxygenase
GENE NAME	PAH
SOURCE ORGANISM	Homo sapiens
TAXONOMY ID	9606 [ <a href="#">NCBI</a> , <a href="#">NEWT</a> ]
LINEAGE	Eukaryota; Metazoa; Chordata; Craniata; Vertebrata; Euteleostomi; Mam

## REFERENCES

[1]	Kwok SCM; Ledley FD; Dilella AG; Robson KJH; Woo SLC <b>Nucleotide sequence of a full-length complementary DNA clone and</b> 1985, <i>Biochemistry</i> , 24, 556-561 <i>Position:</i> SEQUENCE FROM N.A. <i>Comments:</i> tissue=Liver PubMed: <a href="#">2986678</a> ; Medline: <a href="#">85199778</a> ;
-----	---

## COMMENTS

CATALYTIC ACTIVITY	L-phenylalanine + tetrahydrobiopterin + O(2) = L-tyrosine + 4a-hydroxytetrahydrobiopterin
COFACTOR	Ferrous ion.
ENZYME REGULATION	N-terminal region of PAH is thought to contain allosteric binding sites for phenylalanine and to constitute an "inhibitory" regulates the activity of a catalytic domain in the C-terminal portion of the molecule.
PATHWAY	Catabolism of phenylalanine; first (rate-limiting) step.
SUBUNIT	Homodimer.
POLYMORPHISM	The Gln-274 variant occurs on approximately 4% of African-American PAH alleles. The enzyme activity of the variant indistinguishable from that of the wild-type form.
DISEASE	Defects in PAH are the cause of phenylketonuria (PKU) [MIM:261600]. PKU is an autosomal recessive inborn error of phenylalanine metabolism, due to severe phenylalanine hydroxylase deficiency. It is characterized by blood concentration of phenylalanine persistently above 1200 μmol (normal concentration 100 μmol) which usually causes mental retardation, low phenylalanine diet is introduced early in life). They tend to have light pigmentation, rashes similar to eczema, epilepsy hyperactivity, psychotic states and an unpleasant "mousy" odor.
DISEASE	Defects in PAH are the cause of non-phenylketonuria hyperphenylalaninemia (Non-PKU HPA) [MIM:261600]. Non- is a mild form of phenylalanine hydroxylase deficiency characterized by phenylalanine levels persistently below 600 μmol allows normal intellectual and behavioral development without treatment. Non-PKU HPA is usually caused by the combination of a mild hyperphenylalaninemia mutation and a severe one.
DISEASE	Defects in PAH are the cause of hyperphenylalaninemia (HPA) [MIM:261600]. HPA is the mildest form of phenylalanine hydroxylase deficiency.
SIMILARITY	Belongs to the bipterin-dependent aromatic amino acid hydroxylase family.
ONLINE INFORMATION	PAHdb; Phenylalanine hydroxylase locus knowledgebase; Belongs to the bipterin-dependent aromatic amino acid hydroxylase family.

## DATABASE CROSS-REFERENCES

EC	1.14.16.1
EMBL	<a href="#">K03020</a> , AAA60082.1. [ <a href="#">GenBank</a> , <a href="#">DDBJ</a> ] <a href="#">U49897</a> , AAC51772.1. [ <a href="#">GenBank</a> , <a href="#">DDBJ</a> ] <a href="#">S61296</a> , AAD13926.1. [ <a href="#">GenBank</a> , <a href="#">DDBJ</a> ] <a href="#">BC026251</a> , AAH26251.1. [ <a href="#">GenBank</a> , <a href="#">DDBJ</a> ]
GENEW	<a href="#">HGNC:8582</a> , PAH
GO	<a href="#">GO:0004505</a> , F-phenylalanine 4-hydroxylase activity <a href="#">GO:0008652</a> , P-phenylalanine 4-hydroxylase activity
HSC_2DPAGE	<a href="#">P00439</a> , HUMAN
INTERPRO	<a href="#">IPR001273</a> , Aaa_hydroxylase <a href="#">IPR002912</a> , ACT. <a href="#">IPR005961</a> , Phe4hydroxylase
MIM	<a href="#">261600</a>
PDB	<a href="#">1DMW</a> , 2001-03-24. <a href="#">1J8T</a> , 2002-05-22. <a href="#">1J8U</a> , 2002-05-22. <a href="#">1KWQ</a> , 2003-01-28. <a href="#">1LRM</a> , 2002-06-12. <a href="#">1MMK</a> , 2003-09-04. <a href="#">1MMT</a> , 2003-09-04. <a href="#">1PAH</a> , 1999-01-13. <a href="#">2PAH</a> , 1999-10-06. <a href="#">3PAH</a> , 1999-04-27. <a href="#">4PAH</a> , 1999-04-27. <a href="#">5PAH</a> , 1999-04-27.

## KEYWORDS

Oxidoreductase; Monooxygenase; Allosteric enzyme; Phenylketonuria; Phosphorylation; Phenylalanine catabolism; Iron; Disease mutation; Polymorphism; 3D-structure

## FEATURES

Feature	Description	Begin Position	End Position	Length
<a href="#">MOD_RES</a>	Phosphoserine (by PKA) (BY SIMILARITY)	16	16	1
<a href="#">METAL</a>	Iron (BY SIMILARITY)	285	285	1
<a href="#">METAL</a>	Iron (BY SIMILARITY)	290	290	1
<a href="#">METAL</a>	Iron (BY SIMILARITY)	330	330	1
<a href="#">VARIANT</a>	S -> P (in PKU) /FTId=VAR_000869	16	16	1
<a href="#">VARIANT</a>	Q -> L (in HPA) /FTId=VAR_009239	20	20	1
<a href="#">VARIANT</a>	F -> L (in PKU; haplotype 1) /FTId=VAR_000870	39	39	1
<a href="#">VARIANT</a>	(in PKU; haplotypes 9,21) /FTId=VAR_000871	39	39	1
<a href="#">VARIANT</a>	S -> L (in PKU) /FTId=VAR_000872	40	40	1
<a href="#">VARIANT</a>	L -> F (in PKU) /FTId=VAR_000873	41	41	1
<a href="#">VARIANT</a>	L -> P (in PKU; null) /FTId=VAR_009240	41	41	1

## ADDITIONAL INFORMATION FROM iProClass

[Go to iProClass](#)

CROSS-REFERENCE	RefSeq: <a href="#">NP_000268</a> phenylalanine hydroxylase LocusLink: <a href="#">5053</a> phenylalanine hydroxylase(PAH) NCBI GI#: <a href="#">2462722</a> ; <a href="#">189937</a> ; <a href="#">4557819</a>
PIRSF FAMILY	<a href="#">PIRSF000336</a> : phenylalanine 4-mono-oxygenase
GENE ONTOLOGY	<b>Molecular Function</b> <a href="#">GO:0004505</a> : phenylalanine 4-mono-oxygenase activity [ <a href="#">INTERPRO</a> ; evidence: <a href="#">IEA</a> ] [ <a href="#">SPEC</a> ; evidence: <a href="#">IEA</a> ] [PMID:3856322; evidence: <a href="#">TAS</a> ] <a href="#">GO:0016597</a> : amino acid binding [ <a href="#">INTERPRO</a> ; evidence: <a href="#">IEA</a> ] <a href="#">GO:0005506</a> : iron ion binding [ <a href="#">INTERPRO</a> ; evidence: <a href="#">IEA</a> ] <a href="#">GO:0004497</a> : monooxygenase activity [ <a href="#">INTERPRO</a> ; evidence: <a href="#">IEA</a> ] [ <a href="#">SPKW</a> ; evidence: <a href="#">IEA</a> ] <a href="#">GO:0003824</a> : catalytic activity [ <a href="#">SPKW</a> ; evidence: <a href="#">IEA</a> ] <a href="#">GO:0016491</a> : oxidoreductase activity [ <a href="#">SPKW</a> ; evidence: <a href="#">IEA</a> ] <b>Biological Process</b> <a href="#">GO:0009072</a> : aromatic amino acid family metabolism [ <a href="#">INTERPRO</a> ; evidence: <a href="#">IEA</a> ] <a href="#">GO:0008652</a> : amino acid biosynthesis [PMID:3856322; evidence: <a href="#">TAS</a> ] <a href="#">GO:0008152</a> : metabolism [ <a href="#">INTERPRO</a> ; evidence: <a href="#">IEA</a> ] <a href="#">GO:0006559</a> : phenylalanine catabolism [ <a href="#">INTERPRO</a> ; evidence: <a href="#">IEA</a> ] [ <a href="#">SPKW</a> ; evidence: <a href="#">IEA</a> ] [UniProt:P00439; evidence: <a href="#">none</a> ]
ENZYME/FUNCTION	EC 1.14.16.1 <a href="#">EC-IUBMB</a> , <a href="#">KEGG</a> , <a href="#">BRENDA</a> , <a href="#">WIT</a> , <a href="#">MetaCyc</a> <b>Nomenclature:</b> Oxidoreductases; Acting on paired donors, with incorporation or reduction of molecular oxygen, With reduced pteridine as one donor, and incorporation of one atom of oxygen; phenylalanine 4-mono-oxygenase <b>Reaction:</b> L-phenylalanine + tetrahydrobiopterin + O <sub>2</sub> = L-tyrosine + 4a-hydroxytetrahydrobiopterin
PATHWAY	KEGG: Metabolism; Amino Acid Metabolism; Phenylalanine, tyrosine and tryptophan biosynthesis [PATH: <a href="#">hsa00400</a> ].
STRUCTURE	PDB: <a href="#">1PAH</a> (117-424,100.0%) ; <a href="#">1DMW</a> :A(118-424,100.0%) ; <a href="#">More</a> <a href="#">1DMW</a> : <a href="#">SCOP</a> <a href="#">CATH</a> <a href="#">FSSP</a> <a href="#">MMDB</a> <a href="#">PDBsum</a> <a href="#">1IN9</a> : <a href="#">SCOP</a> <a href="#">CATH</a> <a href="#">FSSP</a> <a href="#">MMDB</a> <a href="#">PDBsum</a> <a href="#">More</a>

[http://www.pir.uniprot.org/cgi-bin/upEntry?id=PH4H\\_HUMAN](http://www.pir.uniprot.org/cgi-bin/upEntry?id=PH4H_HUMAN)

caBIG

cancer Biomedical Informatics Grid

# SEED

<b>Developer Center -- POC</b>	U of Chicago – Ed Frank
<b>Adopter Center -- POC</b>	Georgetown – Cathy Wu
<b>caBIG Compliance Level to be Achieved</b>	Silver
<b>Project duration</b>	TBD

# SEED – For subsystem definition

*missing gene in human*

Mycoplasma genitalium [B]	<input type="checkbox"/>												272			
Streptococcus pneumoniae R6 [B]	<input type="checkbox"/>							<a href="#">741</a>	<a href="#">1109</a>	<a href="#">1110</a>	<a href="#">1781</a>	<a href="#">873</a>				
Homo sapiens [E]	<input type="checkbox"/>	<a href="#">18781,24457,3886,421,5247</a>						<a href="#">26330</a>	<a href="#">8649</a>	<a href="#">20554</a>	<a href="#">4329</a>	<a href="#">4329,4427,4428</a>				
Methanocaldococcus jannaschii [A]	<input type="checkbox"/>	<i>Human functional variant</i>														
Haemophilus influenzae Rd KW20 [B]	<input type="checkbox"/>							<a href="#">652</a>	<a href="#">941</a>	<a href="#">601</a>	<a href="#">920</a>	<a href="#">920</a>	<a href="#">621</a>	<a href="#">857</a>		
Aeropyrum pernix [A]	<input type="checkbox"/>	<a href="#">517</a>								<a href="#">1396</a>	<a href="#">1396</a>					
Pyrobaculum aerophilum str. IM2 [A]	<input type="checkbox"/>	<a href="#">2402</a>								<a href="#">831</a>	<a href="#">831</a>	<a href="#">608</a>				
Buchnera aphidicola str. APS (Acyrthosiphon pisum) [B]	<input type="checkbox"/>	<a href="#">196</a>							<a href="#">566</a>	<a href="#">195</a>					<a href="#">554</a>	<a href="#">202</a>
Synechocystis sp. PCC 6803 [B]	<input type="checkbox"/>	<a href="#">2844</a>							<a href="#">2956</a>	<a href="#">1587</a>	<a href="#">1768</a>	<a href="#">1542</a>	<a href="#">1366</a>	<a href="#">1366</a>	<a href="#">2569</a>	<a href="#">2904</a>
Organism	Variant Code	ASPDC	KPHMT	KPRED	KARED	PBAL	PANF		PANK	PPCS	PPCDC		PPAT	DPCK		
Saccharomyces cerevisiae [E]	<input type="checkbox"/>		<a href="#">3766</a>	<a href="#">1690</a>	<a href="#">3279</a>	<a href="#">5013</a>				<a href="#">2331</a>	<a href="#">5073</a>	<a href="#">2556,5673,5824</a>	<a href="#">4971</a>	<a href="#">2011</a>		
Bacillus anthracis str. A2012 [B]	<input type="checkbox"/>	<a href="#">2391</a>		<a href="#">2694</a>	<a href="#">2247,2663</a>	<a href="#">4407,4496</a>			<a href="#">3725,963</a>	<a href="#">4786</a>	<a href="#">4786</a>		<a href="#">4914</a>	<a href="#">5560</a>		
Helicobacter pylori J99 [B]	<input type="checkbox"/>	<a href="#">30</a>	<a href="#">364</a>		<a href="#">310</a>	<a href="#">6</a>				<a href="#">792</a>	<a href="#">775</a>	<a href="#">775</a>	<a href="#">1364</a>	<a href="#">766</a>		
Thermotoga maritima [B]	<input type="checkbox"/>	<a href="#">931</a>	<a href="#">1710</a>		<a href="#">544</a>	<a href="#">1067</a>				<a href="#">875</a>	<a href="#">1671</a>	<a href="#">1671</a>	<a href="#">734</a>	<a href="#">1372</a>		
Mycobacterium tuberculosis H37Rv [B]	<input type="checkbox"/>	<a href="#">3604</a>	<a href="#">2227</a>	<a href="#">2575</a>		<a href="#">3605</a>				<a href="#">1094,3603</a>	<a href="#">1393</a>	<a href="#">1393</a>	<a href="#">2967</a>	<a href="#">1633</a>		
Bacillus subtilis subsp. subtilis str. 168 [B]	<input type="checkbox"/>	<a href="#">2245</a>	<a href="#">2247</a>	<a href="#">1446</a>	<a href="#">2832</a>	<a href="#">2246</a>	<a href="#">323,667</a>				<a href="#">2381</a>	<a href="#">1572</a>	<a href="#">1572</a>	<a href="#">1504</a>	<a href="#">2909</a>	
Escherichia coli K12 [B]	<input type="checkbox"/>	<a href="#">131</a>	<a href="#">134</a>	<a href="#">422</a>	<a href="#">3704</a>	<a href="#">133</a>	<a href="#">3200</a>				<a href="#">3890</a>	<a href="#">3575</a>	<a href="#">3575</a>	<a href="#">3570</a>	<a href="#">103</a>	

*Human functional variant*

# **SIG 3:**

## **Microarray Repositories**

### **(3, 47, 15)**

The mission of the Microarray Repositories SIG is to identify and prioritize the needs of the larger cancer research community with respect to the capture, storage, and utilization of microarray data and related types of genetic data. Specifically, we will address the need for:

- A database for the storage and retrieval of microarray data and related data types that can be incorporated into the larger scheme of federated databases that store clinically relevant data and other relevant data types.
- Software that facilitates the capture of important microarray experimental information and automatically loads it into a database.
- Software that facilitates the querying of microarray databases and the retrieval of data.
- Consumers and producers of microarray data to readily exchange data

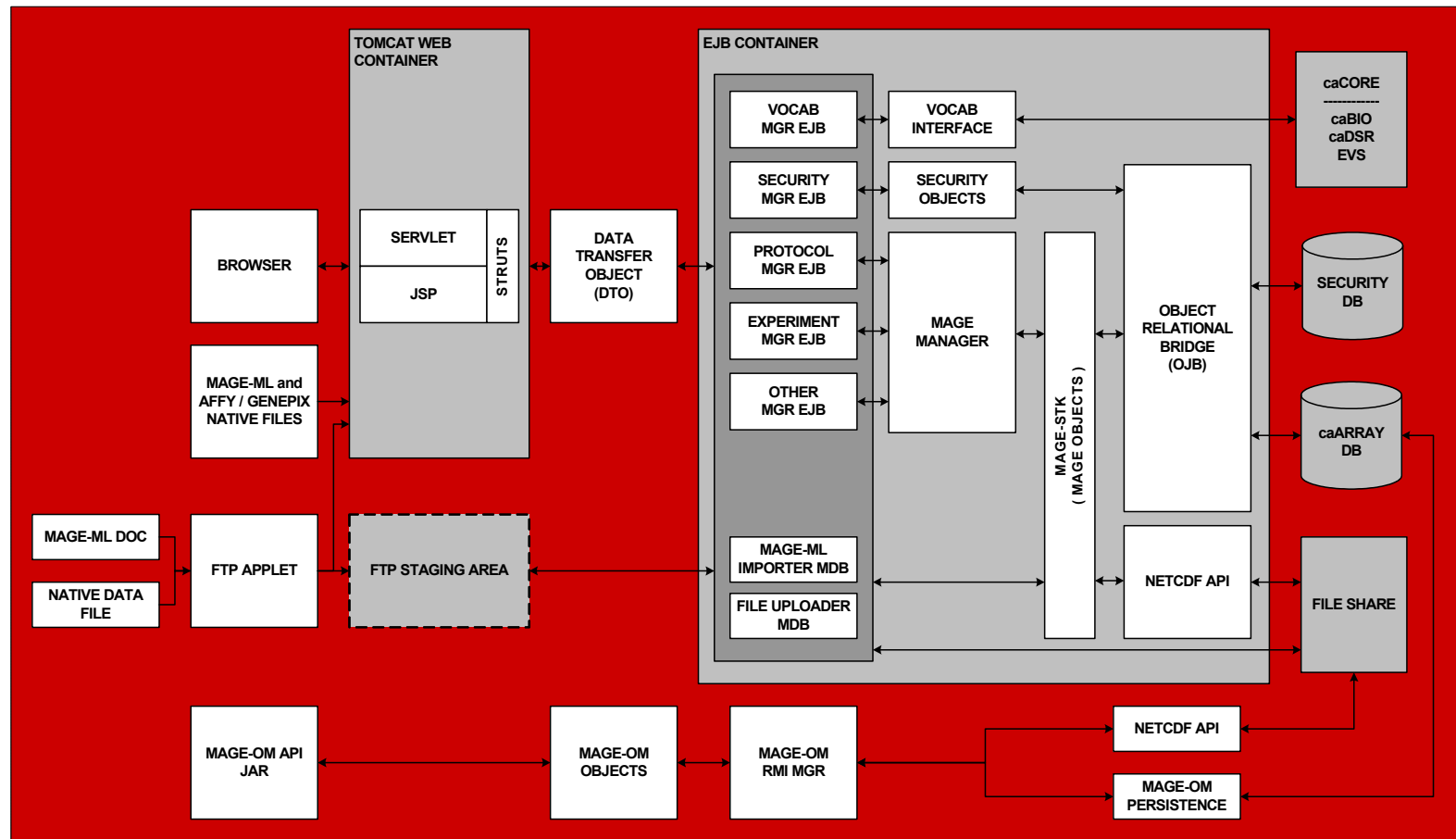
## caArray

<b>Developer Center -- POC</b>	NCICB – Mervi Heiskanen
<b>Adopter Center -- POC</b>	Wistar – Louise Showe NYU – Judith Goldberg Georgetown – Arnie Miles
<b>caBIG Compliance Level to be Achieved</b>	N/A
<b>Project duration</b>	7 months

# caArray Project Objectives

- ▶ Create a Test Procedures document
- ▶ Install the caArray software and create an Installation Guide
- ▶ Test the experiment annotation, data loading and data transfer functionality of caArray using test data provided by NCICB and local data
- ▶ Review and update the Usage Guide provided by NCICB
- ▶ Train users; Review and update the training materials provided by NCICB

# caArray Architecture





## NCI-60 Data Sharing

<b>Developer Center -- POC</b>	NCI-CCR – David Kane
<b>Adopter Center -- POC</b>	MSKCC – Alex Lash
<b>caBIG Compliance Level to be Achieved</b>	Silver
<b>Project duration</b>	12 months

# NCI-60 Project Objectives

- ▶ Create a Work Plan for the project
- ▶ Create a Risk Management Matrix for the project
- ▶ Document a Test Approach that ensures that data are transformed accurately
- ▶ Create the following data resources
  - XML-based representation of NCI-60 cell line characterization
  - MAGE-ML files and source data files for Affymetrix U95 and U133 NCI-60 datasets
  - MAGE-ML toxicology files for two NCI-60 drug screen datasets
  - MAGE-ML files and source data files for the NCI-60 9,706 clone cDNA gene expression data sets
  - SKYWEB files for the NCI-60 karyotyping data
  - MAGE-ML files and source data files for the aCGH NCI-60 datasets
- ▶ As appropriate and feasible, upload these files to the current caBIG reference implementation for these standards, e.g. caArray.
- ▶ Execute on Test Approach

## Zebrafish Microarray Data Sharing

<b>Developer Center -- POC</b>	Thomas Jefferson – Jack London
<b>Adopter Center -- POC</b>	MSKCC – Alex Lash
<b>caBIG Compliance Level to be Achieved</b>	Silver
<b>Project duration</b>	TBD

## **SIG 4: Pathways (3, 35, 10)**

The Pathway SIG endeavors to support basic research and ICR tool development by helping to provide the cancer research community with easy access to pathway data and commonly used pathway analysis tools.

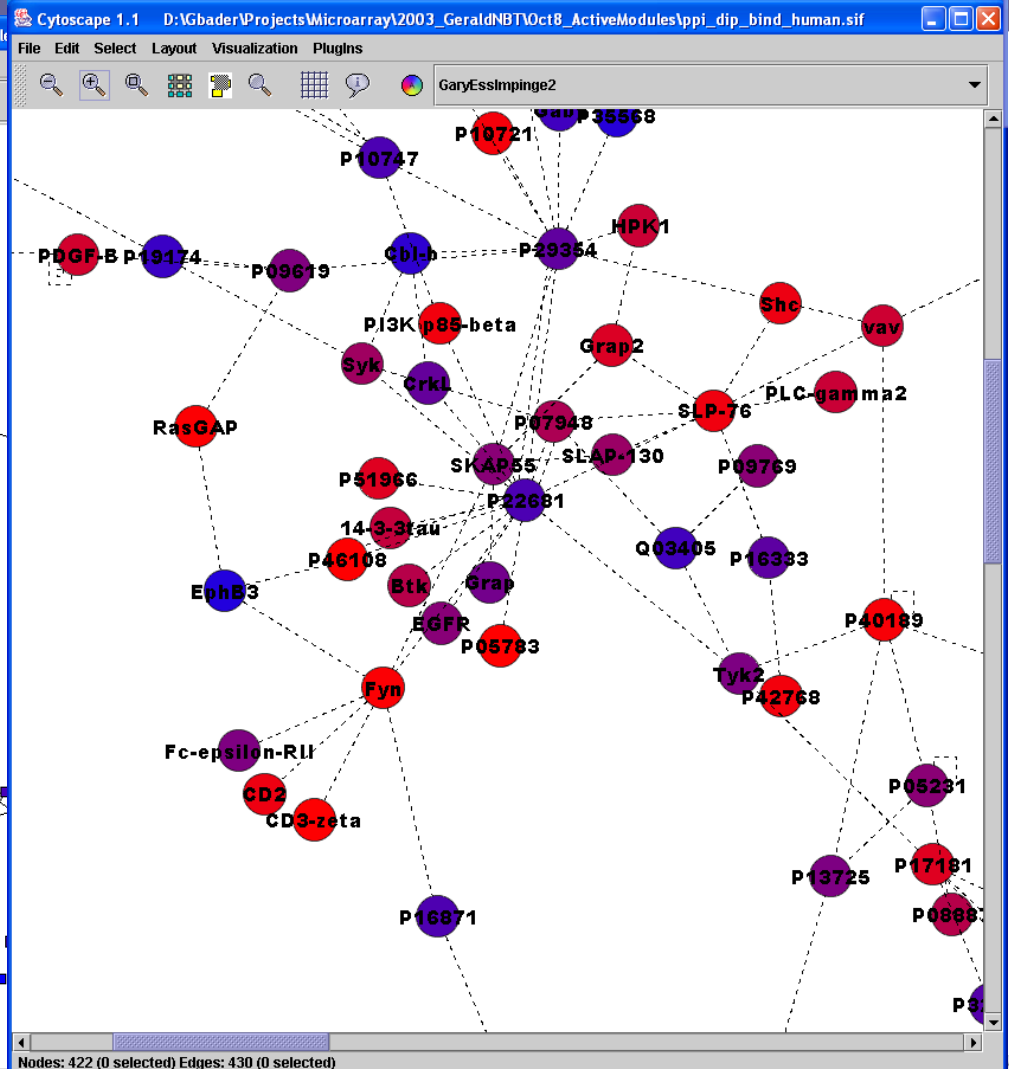
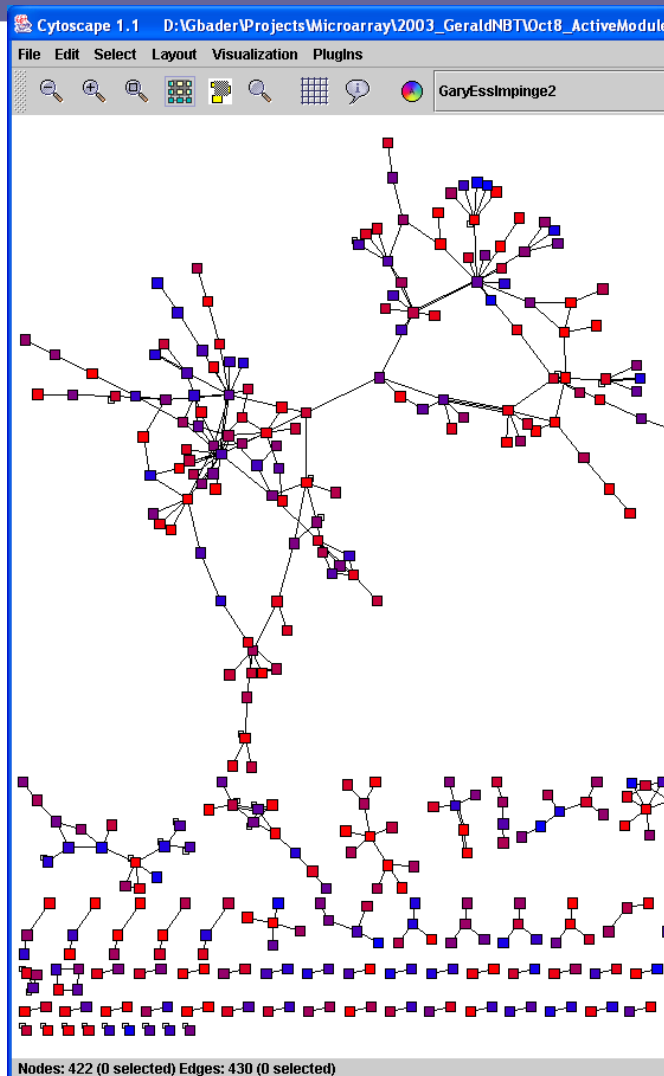
## Pathways Tool Development

<b>Developer Center -- POC</b>	MSKCC – Gary Bader
<b>Adopter Center -- POC</b>	OHSU – Shannon McWeeney
<b>caBIG Compliance Level to be Achieved</b>	Silver
<b>Project duration</b>	12 months

# Pathways Project Objectives

- ▶ Develop a Functional Requirements and Design Specification document, in collaboration with the Adopter Center and the cross-cutting Workspaces
- ▶ Create a Risk Management Matrix for the project
- ▶ Document a Test Approach that ensures requirements are met
- ▶ Write code to achieve the following milestones:
  - Build I/O functionality for BioPAX format in the Data Services Framework used by cPath and Cytoscape
  - Expand cPath data model to BioPAX 1.0
  - Update Cytoscape to show a simple view of BioPAX data
  - Create Cytoscape plugin for loading caBIG gene annotation and expression data
- ▶ Execute on Test Approach

# Cytoscape Screen Shots



## Quantitative Analysis of Pathways in Cancer (QPACA)

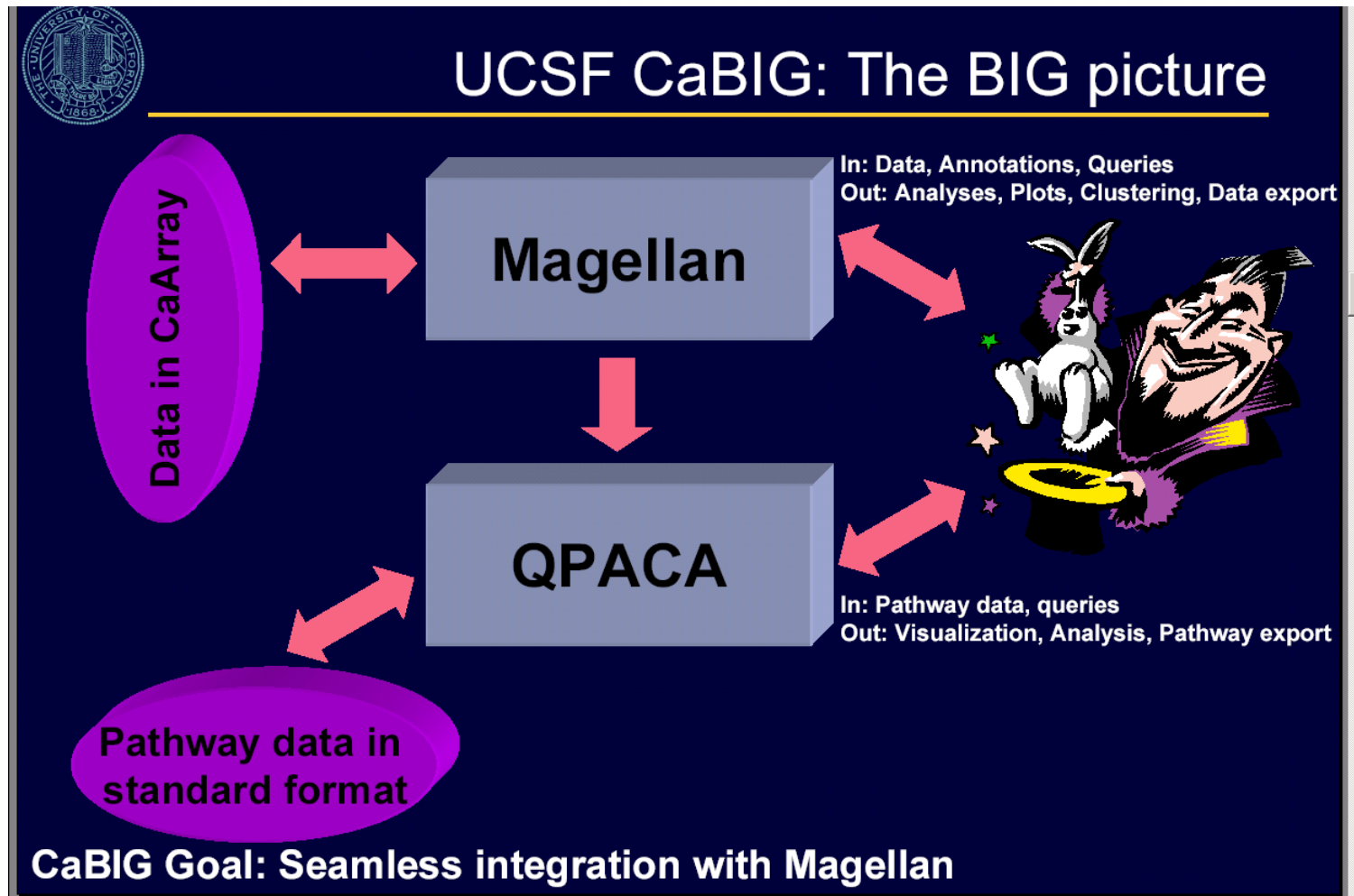
<b>Developer Center -- POC</b>	MSKCC – Gary Bader
<b>Adopter Center -- POC</b>	OHSU – Shannon McWeeney
<b>caBIG Compliance Level to be Achieved</b>	Silver
<b>Project duration</b>	12 months



# QPACA Project Objectives

- ▶ Develop a Functional Requirements and Design Specification document, in collaboration with the Adopter Center(s) and the cross-cutting Workspaces
- ▶ Create a Risk Management Matrix for the project
- ▶ Document a Test Approach that ensures requirements are met
- ▶ Make current version of QPACA available to the Adopter for evaluation
- ▶ Write code to achieve the following milestones:
  - Make QPACA interoperable with Magellan
  - Refine QPACA statistical methods
  - Make QPACA interoperable with caBIG-defined Pathways Exchange Standards
- ▶ Execute on Test Approach

# QPACA and caBIG



## Reactome Project

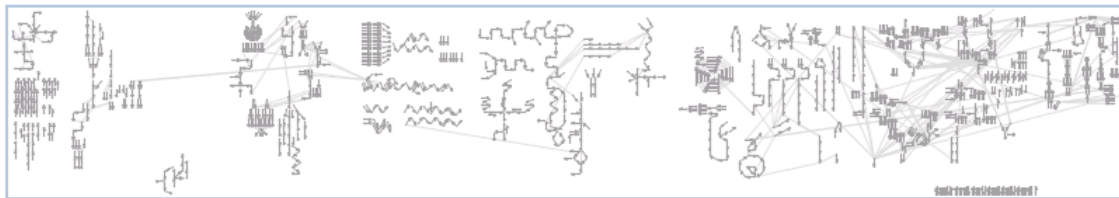
<b>Developer Center -- POC</b>	Cold Spring Harbor – Brian Gilman
<b>Adopter Center -- POC</b>	MSKCC – Gary Bader
<b>caBIG Compliance Level to be Achieved</b>	Silver
<b>Project duration</b>	TBD

# Reactome Home Page


[About](#)
[TOC](#)
[User Guide](#)
[Data Model](#)
[Schema](#)
[Extended search](#)
[Pathfinder](#)
[Download](#)
[Linking](#)
[Citing](#)

Find  with  in

## Reactome - a knowledgebase of biological processes



<b>Apoptosis</b> <a href="#">Hsa</a>   <a href="#">Mmu</a>   <a href="#">Rno</a>   <a href="#">Dre</a>   <a href="#">Fru</a>   <a href="#">Gga</a>	<b>Cell Cycle Checkpoints</b> <a href="#">Hsa</a>   <a href="#">Mmu</a>   <a href="#">Rno</a>   <a href="#">Dre</a>   <a href="#">Fru</a>   <a href="#">Gga</a>	<b>Cell Cycle, Mitotic</b> <a href="#">Hsa</a>   <a href="#">Mmu</a>   <a href="#">Rno</a>   <a href="#">Dre</a>   <a href="#">Fru</a>   <a href="#">Gga</a>	<b>DNA Repair</b> <a href="#">Hsa</a>   <a href="#">Mmu</a>   <a href="#">Rno</a>   <a href="#">Dre</a>   <a href="#">Fru</a>   <a href="#">Gga</a>
<b>DNA Replication</b> <a href="#">Hsa</a>   <a href="#">Mmu</a>   <a href="#">Rno</a>   <a href="#">Dre</a>   <a href="#">Fru</a>   <a href="#">Gga</a>	<b>Gene Expression</b> <a href="#">Hsa</a>   <a href="#">Mmu</a>   <a href="#">Rno</a>   <a href="#">Dre</a>   <a href="#">Fru</a>   <a href="#">Gga</a>	<b>Hemostasis</b> <a href="#">Hsa</a>   <a href="#">Mmu</a>   <a href="#">Rno</a>   <a href="#">Dre</a>   <a href="#">Fru</a>   <a href="#">Gga</a>	<b>Insulin receptor mediated signalling</b> <a href="#">Hsa</a>   <a href="#">Mmu</a>   <a href="#">Rno</a>   <a href="#">Dre</a>   <a href="#">Fru</a>
<b>Lipid metabolism</b> <a href="#">Hsa</a>   <a href="#">Mmu</a>   <a href="#">Rno</a>   <a href="#">Dre</a>   <a href="#">Fru</a>   <a href="#">Gga</a>	<b>Metabolism of amino acids and related nitrogen-containing molecules</b> <a href="#">Hsa</a>   <a href="#">Mmu</a>   <a href="#">Rno</a>   <a href="#">Dre</a>   <a href="#">Fru</a>   <a href="#">Gga</a>	<b>Metabolism of glucose, other sugars, and ethanol</b> <a href="#">Hsa</a>   <a href="#">Mmu</a>   <a href="#">Rno</a>   <a href="#">Dre</a>   <a href="#">Fru</a>   <a href="#">Gga</a>	<b>mRNA Processing</b> <a href="#">Hsa</a>   <a href="#">Mmu</a>   <a href="#">Rno</a>   <a href="#">Dre</a>   <a href="#">Fru</a>   <a href="#">Gga</a>
<b>Nucleotide metabolism</b> <a href="#">Hsa</a>   <a href="#">Mmu</a>   <a href="#">Rno</a>   <a href="#">Dre</a>   <a href="#">Fru</a>   <a href="#">Gga</a>	<b>Oxidative decarboxylation of pyruvate and TCA cycle</b> <a href="#">Hsa</a>   <a href="#">Mmu</a>   <a href="#">Rno</a>   <a href="#">Dre</a>   <a href="#">Fru</a>   <a href="#">Gga</a>	<b>Transcription</b> <a href="#">Hsa</a>   <a href="#">Mmu</a>   <a href="#">Rno</a>   <a href="#">Dre</a>   <a href="#">Fru</a>   <a href="#">Gga</a>	<b>Translation</b> <a href="#">Hsa</a>   <a href="#">Mmu</a>   <a href="#">Rno</a>   <a href="#">Dre</a>   <a href="#">Fru</a>   <a href="#">Gga</a>

### About Reactome



The **Reactome** project is a collaboration among Cold Spring Harbor Laboratory, The European Bioinformatics Institute, and The Gene Ontology Consortium to develop a curated resource of core pathways and reactions in human biology. The information in this database is authored by biological researchers with expertise in their field, maintained by the Reactome editorial staff, and cross-referenced with with PubMed, GO, and the sequence databases Ensembl and UniProt.

Reactome is a free on-line resource, and Reactome software is open-source. However, please take note of our [disclaimer](#).

[More...](#)

### News and Notes

- October 27, 2004: 11th Release of Reactome**  
 New modules released today are Apoptosis and Hemostasis.
- October 12, 2004: Reactome now downloadable in SBML format**  
 Human reactions in Reactome can now be downloaded in **SBML** format. Also, SBML version of Reactome events is available at the bottom of each event page.
- [More...](#)

The development of Reactome is supported by grant R01 HG002639 from the National Human Genome Research Institute at the US National Institutes of Health, grant LSHG-CT-2003-503269 from European Union (6th Framework Programme) and a subcontract from the NIH-funded Cell Migration Consortium.

© 2003 Cold Spring Harbor Laboratory and European Bioinformatics Institute. All rights reserved.

## **SIG 5/6: Proteomics (3, 43, 10)**

The Proteomics Tools SIG is focused on:

- Tools and technologies which are necessary for cancer centers to store, annotate, and analyze the growing proteomics data sets
- Providing flexible tools for data and metadata storage, so that emerging technologies can be incorporated into developed systems
- Integrating proteomics data with other data through use of appropriate CDEs and architectures

## Proteomics LIMS

<b>Developer Center -- POC</b>	Fox Chase – Michael Ochs
<b>Adopter Center -- POC</b>	Moffitt – Steven Eschrich
<b>caBIG Compliance Level to be Achieved</b>	Silver
<b>Project duration</b>	12 months

# Proteomics LIMS Project Objectives

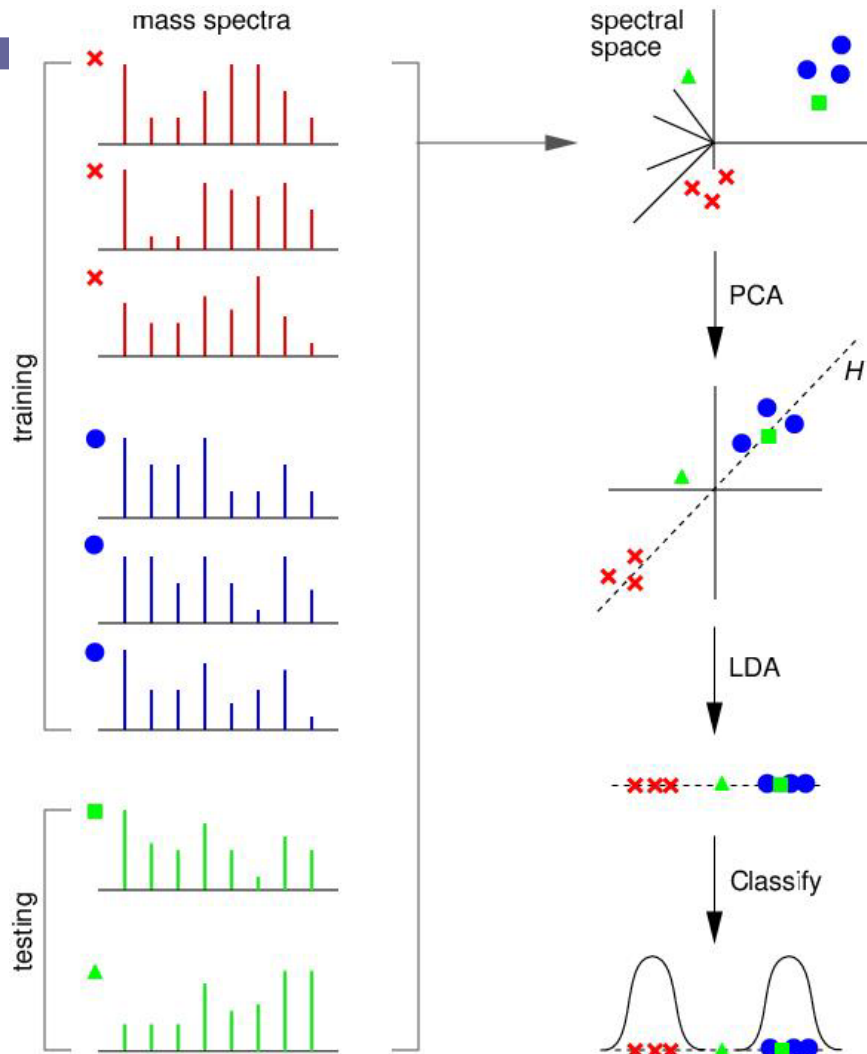
- ▶ Develop use cases to represent the lab processes around 2D gel electrophoresis
- ▶ Develop a Functional Requirements and Design Specification document, in collaboration with the Adopter Center and the cross-cutting Workspaces
- ▶ Create a UML model to represent the system
- ▶ Create a Risk Management Matrix for the project
- ▶ Document a Test Approach that ensures requirements are met
- ▶ Create mock-ups of the system user interfaces for review by the Adopter
- ▶ Execute on Test Approach
- ▶ Deploy prototype system to the Developer and Adopter sites

## Q5

<b>Developer Center -- POC</b>	Dartmouth – David Jewell
<b>Adopter Center -- POC</b>	OHSU – Shannon McWeeney
<b>caBIG Compliance Level to be Achieved</b>	Silver
<b>Project duration</b>	TBD



## Q5: Disease Classification by Mass Spec Pattern Recognition



### Proteome Analysis

Expression analysis: proteins (mass spectrometry)

Algorithm uses PCA followed by LDA

probabilistic classification of healthy vs. disease whole serum samples using mass spectrometry

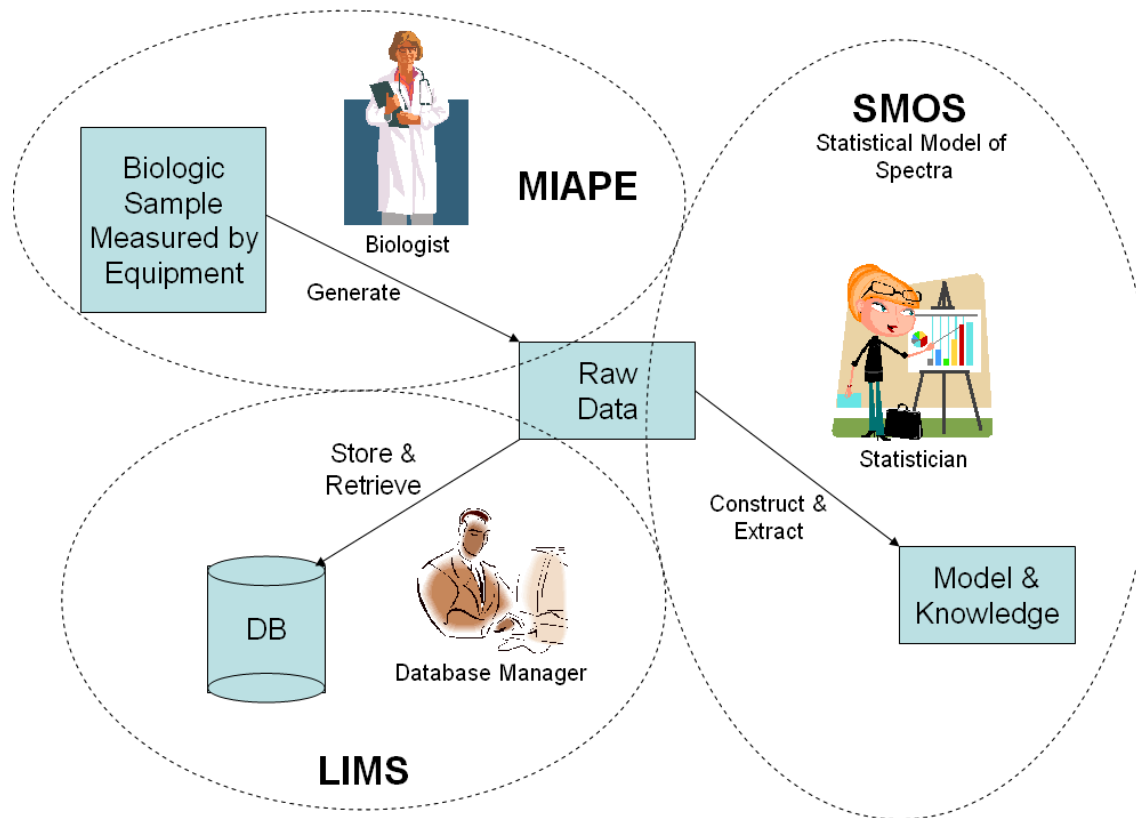
## RProteomics

<b>Developer Center -- POC</b>	Duke – Patrick McConnell
<b>Adopter Center -- POC</b>	OHSU – Shannon McWeeney Penn – David Fenstermacher
<b>caBIG Compliance Level to be Achieved</b>	Gold
<b>Project duration</b>	9 months

# RProteomics Project Objectives

- ▶ In collaboration with the Adopter centers, develop use cases for RProteomics
- ▶ Develop Functional Requirements and Specification document, in collaboration with the Adopter Centers and the Cross-Cutting Workspaces
- ▶ Create a Risk Management Matrix for the project
- ▶ Document a Test Approach that ensures requirements are met
- ▶ Describe the RProteomics data structures in UML and as common data elements
- ▶ Create source code for
  - R and Java data structures,
  - Java-to-R bridge
  - Java and R grid wrappers
  - Client application
- ▶ Document best practices for integrating Java and R
- ▶ Present project to ICR and Architecture Workspaces

# Relationship between MIAPE and SMOS



## **SIG 6/6: Translational (1, 40, 10)**

The Translational Tools Special Interest Group is focused on:

- Tools and technologies which are necessary for cancer centers to integrate clinical data with experimental data, as well as experimental design tools that provide assistance to biomedical investigators embarking on translational research methodology
- Assuring that clinical data and studies more effectively utilize genomics and proteomics research data in cancer research and patient care
- Creating guidelines to aid in the design of experimental studies

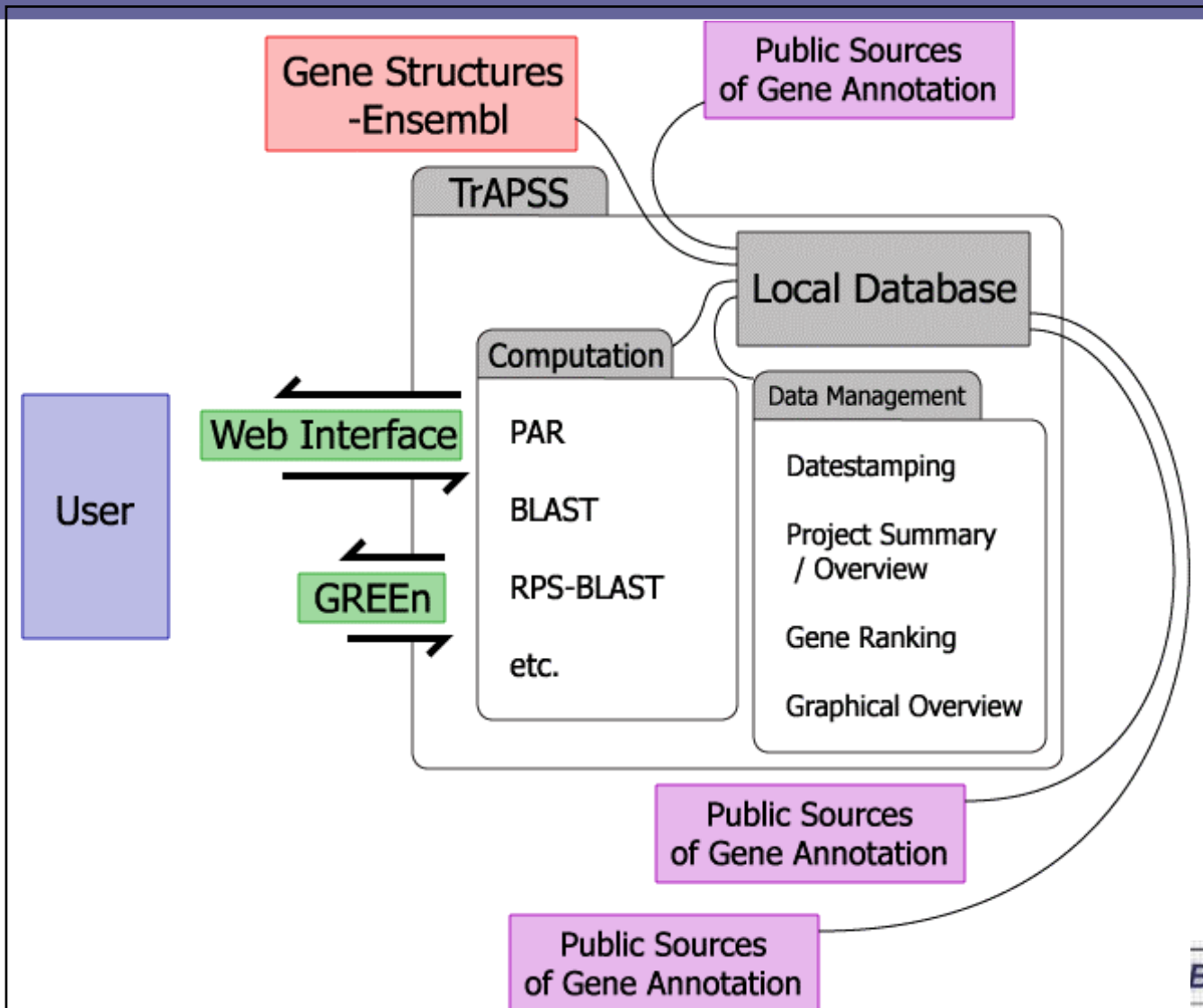
# Transcript Annotation Prioritization and Screening System (TrAPSS)

<b>Developer Center -- POC</b>	U of Iowa – Terry Braun
<b>Adopter Center -- POC</b>	Wistar – Harold Riethman
<b>caBIG Compliance Level to be Achieved</b>	Silver
<b>Project duration</b>	12 months

# TrAPSS Project Objectives

- ▶ Create a Risk Management Matrix for the project
- ▶ Make current version of TrAPSS available to the Adopter and train them on the use of the system
- ▶ Develop a Use Case document, in collaboration with the Adopter Center
- ▶ Develop a Functional Requirements and Specification Document in collaboration with the Adopter Center and the cross-cutting Workspaces
- ▶ Document a Test Approach that ensures requirements are met
- ▶ Document and implement necessary changes to the TrAPSS database
- ▶ Implement the changes necessary to the TrAPSS modules as defined in the Requirements and Specification Document
- ▶ Execute on Test Approach
- ▶ Deploy system to Adopter site
- ▶ Create a Developer's guide

# TrAPSS System





# ICR Workspace at a Glance

## - Conclusion::Goals Defined (drafted)

### ► One Year Goals:

- The majority of projects will target Silver Level compliance, as outlined in the caBIG Compatibility Guidelines document. A few projects will be selected as reference implementations of the Grid architecture. All projects will comply with the caBIG principles of open source, open access, open development and federation.

### ► Three year goals:

- 1. Bring more selected data sources and applications into Gold-level compliance, i.e. make them grid resources.
- 2. Show how a federated set of resources that are syntactically and semantically compatible can be used to perform powerful analyses across the cancer research community.

### ► Five year Goals

- Connect the ICR data and applications to resources in the Tissue Banks and Pathology Tools Workspace as well as the Clinical Trials Management System Workspace

# ICR Workspace at a Glance

## - Conclusion::Perceived Challenges

- ▶ “Tightening” the group – how do projects “fit” together?
- ▶ “Cementing” the Developer-Adopter relationships (More use-cases need to be developed jointly)
- ▶ Increasing awareness crossing...
  - Workspaces (beyond liaisons)
  - SIGs (distinction?/commonality?)
- ▶ Uniform training
- ▶ Uniform adoption of caBIG Architecture, CDEs and practices
- ▶ Balancing Software Engineering with Science
  - Discussions
  - Effort
- ▶ A “Vision” for the “Grid”